

- Future Activities
- Programming Environment
- Performance of Level 3 routines
- Overview of SciDAC software activities

June 23th 2004, Fermilab
For **QCD**QCDOC** collaboration
(Brookhaven National Laboratory/Columbia University)
Chulwoo Jung**

The Status of User Software on QCDOC

QCDOC Collaboration

Columbia University: N. Christ, R. Ma whinney, A. Yamaguchi,

H. W. Lin, G. F. Liu, X. Liao, C. Kim,

L. Levkova, S. Cohen, S. Li

P. Boyle, B. Joo, M. Clark

C. Jung, K. Petrov

T. Wetzig (Regensburg), S. Ohta (RIKEN)

A. Garra, D. Chen

IBM Yorktown:

RBC:

BNL:

UKQCD:

- QCDOC should be able to sustain performance when the number of processors increases to several thousands.
 - Use **On-Chip Technology** to combine computation and communication units in one chip (**QCDOC = QCD On a Chip**)
 - ASIC (Application-Specific Integrated Circuit) combines existing IBM components and **QCDOC**.
 - **specific, custom-designed logic:**
 - 500 MHz PowerPC 440 processor core with 64-bit, 1 GFlops FPU (32+32KB L1 cache)
 - 4 MB on-chip memory (embedded DRAM), accessed through custom-designed prefetching
 - eDRAM controller (PEC)
 - 6-dimensional nearest-neighbor, low latency serial communications unit (**SCU**) with aggregate bandwidth of 12 Gbit/s with software partitioning capability for physics communications.
 - Auxiliary fast ethernet network for booting, IO, and diagnostics.
 - ~ 1 US-\$ per sustained MFlops, Scalable to several 10 TFlops
- (Planevy talk by P. Boyle)

Hardware overview of QCDOC

Software objective: Create a unified programming environment that will enable US lattice computers to achieve high efficiency on diverse multi-teraflop scale hardware.

Existing codes for lattice QCD (CPS, MILC, Chroma/SZIN,) can be run on various hardware by adopting to programming environment provided by SciDAC Software effort.

Participant list

Arizona	Doug Toussaint	Illinois	Celso Mendez	Eric Gregory	Illinois	Daniel Reed	Robert Edwards	Chip Watson	Walt Akers	Jie Chen	Konstantin Petrov	Jlab	Chulwoo Jung	BNL	John Negel	MIT	Andrew Pochinsky	Carleton De Tar	Jim Simone	Don Holmgren	Bob Mawhinney	MIT	James Osborn	Amitoj Singh	Fermilab	
Boston U.				Hartmut Neff	Jlab																					
Boston U.																										
BNL																										
Columbia																										
Fermilab																										
Fermilab																										

SciDAC Software Program

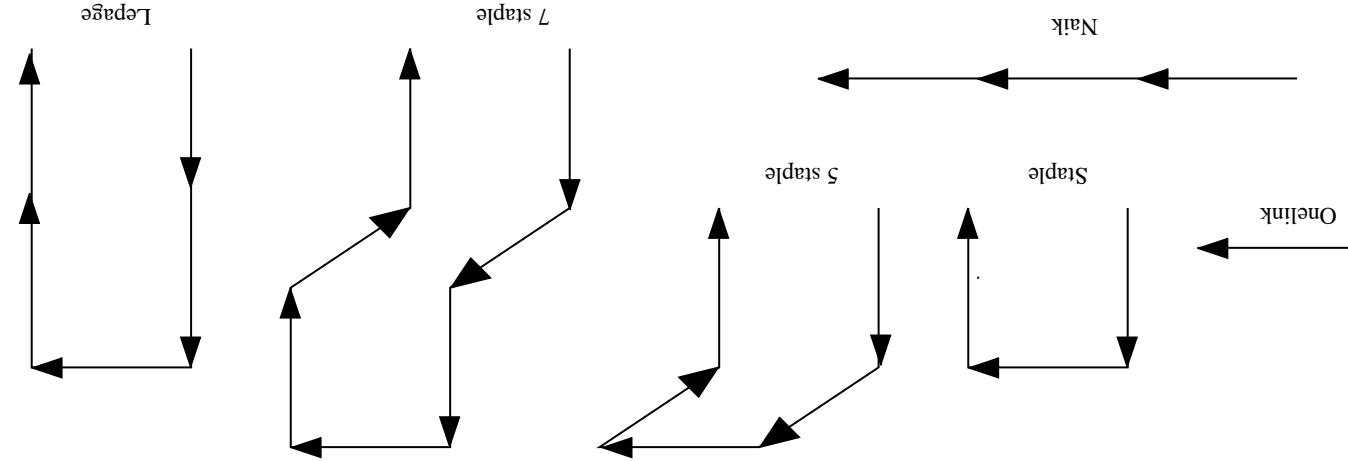
- **Level 3:** Highly optimized routines for Lattice QCD (**QOP**)
 - Wilson, Clover, DWF (P. Boyle)
 - Staggered (C. Christian)
 - Asqtad (improved staggered) (CJ, H. Lin)
 - Asqtad force term (CJ, Z. Sroczynski)
 - Asqtad inverter has been tested with MILC code (CJ, E. Gregory).
 - Level 3 interface is being implemented for Asqtad inverter/force term (CJ).
 - Asqtad HMC is being tested and also being used for hardware testing.
 - Wilson Slash is also integrated into Chroma/QDP++(P. Boyle, B. Joo).
- Currently available in **Columbia Physics System (CPS)**:

QCD-API on QCDOC

- Level 2:

- QCD Data Parallel QDP/QDP++: Lattice-wide operations with both communication and computation (J. Osborn, B. Joo, R. Edwards)
- example: parallel transport $\chi(x) = U^u(x)\chi(x + \ell)$
- Both QDP(C) and QDP++(C++) has been tested with QMP.
- QIO: Lattice data IO (B. Joo, K. Petrow)
- Level 1:
- QMP: Internode communication (CJ)
- QLA: Single node linear Algebra (J. Osborn).

Asqtad action is one of the improved staggered fermion action, with smaller lattice spacing error and flavor mixing.



Computation counts for various terms needed for HMC (per site):

Dirac inverter: $\sim (1150 \times (n \sim 1000))$

Fermion Force term: $\sim 434,000$ for 2+1 flavor. Calculation of the force term mostly consists of parallel transport and $SU(3)$ vector outer products.

Farlink: $\sim 60,000$ Gauge Force: $\sim 150,000$

Vector outer product: $P(x) = X(x)X^\dagger(x)$

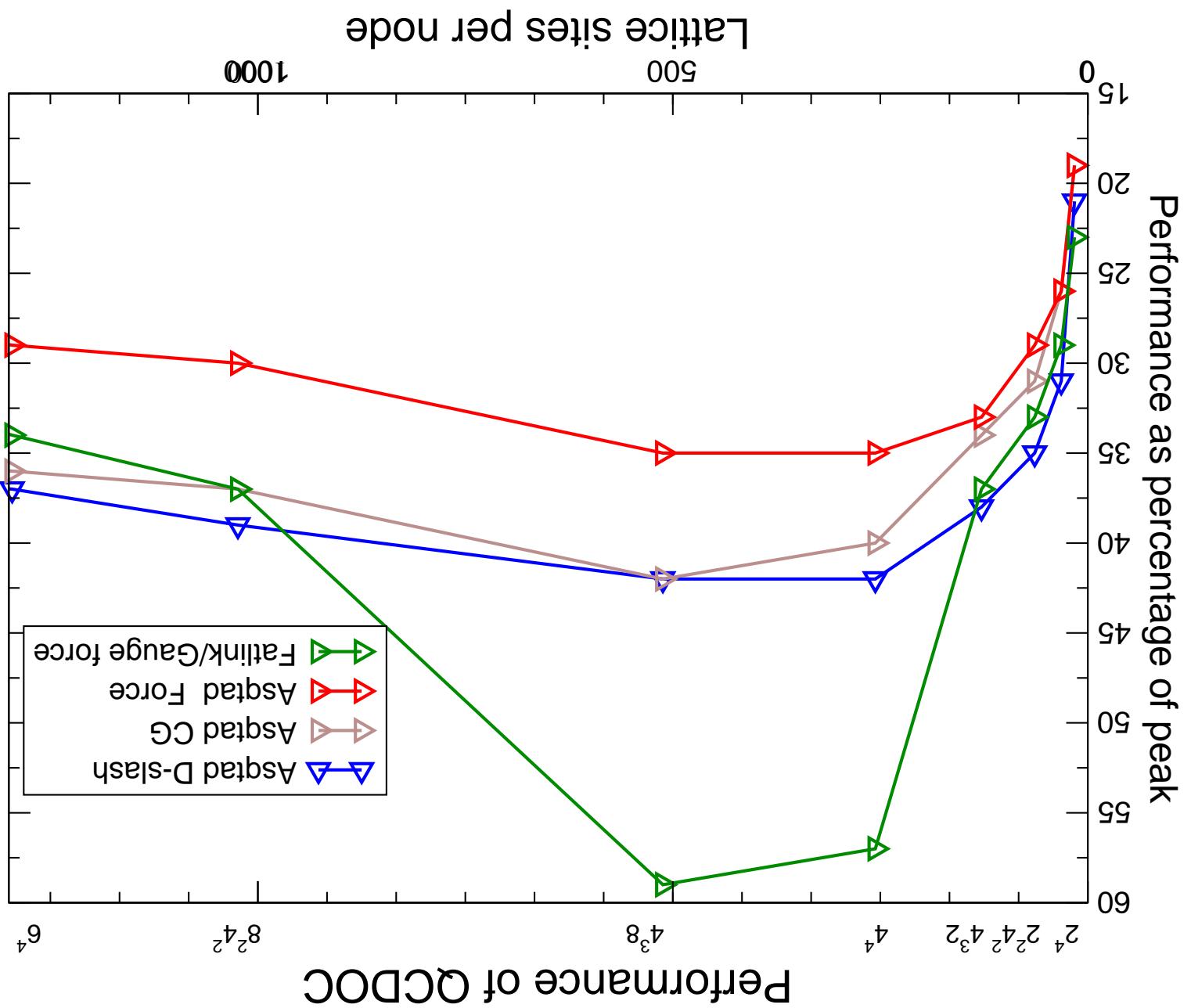
Parallel transport: $X_\mu(x) = U_\mu(x)X(x + h)$

Farlink : $\sim 60,000$ Gauge Force: $\sim 150,000$

(*): multi-dimensional SU(3) parallel transport

Precision	Sites/node	Optimized	MILC	MILC/QDP
Application		D	Double	Single
Machine size	Sim	128	128	16
Asqtad Dirac/CG	2 ⁴ 4 ⁴ 6 ⁴ 8 ⁴	20% 19% 40% 42%	9.5% 19% 36% 37%	15% 15.6% 28% 29%
Machine size		128	128	16
Asqtad Force	2 ⁴ 4 ⁴ 6 ⁴ 8 ⁴	19% 35% 29% 20%	10.0% 10.8% 10.8% 8.2%	
Asqtad Gauge/	2 ⁴ 4 ⁴ 6 ⁴ 8 ⁴	23% 57% 34% 7.2%	8.1% 8.1% 8.1% 7.2%	
Asqtad Force*	2 ⁴ 4 ⁴ 6 ⁴ 8 ⁴		27% 27% 27% 27%	8.5% 8.5% 8.5% 8.5%

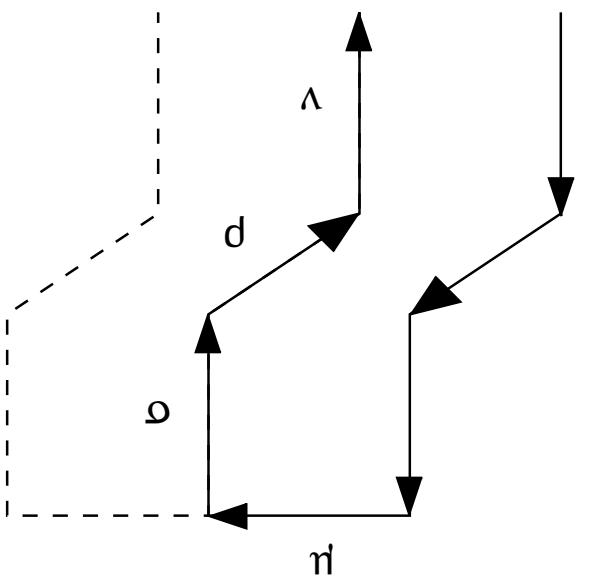
Performance of QCDOC



- For most computation routines, memory bandwidth is bottleneck (EDRAM/DDR → L1 cache → CPU). Changing to single precision only improves memory → L1 bandwidth.
- Sequential read access maximizes memory bandwidth. (Scatter don't gather!)
- Assembly routines are written by C++ programs (CJ, P. Boyle(BAGEL))
- For asqtad dirac, computation is divided into local/non-local. (Usual forward-backward split increases the amount of fermion fields communication)
- Start communication
 - Local computations
 - Wait for communication
 - Non-local computations

Optimization strategy for QCDOC

- For the force term and fatlink generation, multi-dimensional parallel transport is used to improves communication bandwidth. This should be balanced with memory usage.
- For example, 7 plaquette with $+x, -y, -t, +z$ can be calculated simultaneously with $+x, -y, -t, -z$, etc. Also, you can save on number of computations by calculating $-x, +y, +t, -z$, etc. (Strategy for QDP?)
- Using multi-directional parallel transport whenever possible is a fast way to generate reason-ably optimized code for small local volume. (Strategy for QDP?)



- QoS - Custom-designed Operating system (P. Boyle)
 - Boot (Ethernet/JTAG)
 - Load/Run Apps (Ethernet)
 - Hardware diagnostics
 - Command line argument passing
 - File I/O - NFS heavily tested (currently to the front end)
 - Shell support (qsch : tcsh + QCDQC-specific commands) (K. Petrov/ P. Boyle)
 - Profiling, Source level debugging
- Standard GNU tools (Compiler, Linker, Binutils...) at the front end
 - Multi-processor debugger - Riscwatch
 - High-performance compilers - XLC
- Commercial PowerPC tools
 - Multi-processor debugger - Riscwatch
 - High-performance compilers - XLC

Programming Environment

- The goal of MPI is to provide portable, low-latency, high-bandwidth communication routines suitable for Lattice QCD.
- Point-to-point communication
 - Nonblocking
 - Simultaneous, multidirectional transfer capability
 - Chained block/stripped transfer capability
 - Separate routes for the initialization and commencement of transfers → opened communication channels can be reused, minimize overhead for repeated transfers.
 - Global operation
 - Global reduction
 - General binary reduction available
 - {Sum, Max, Min} {Single precision, Double precision} {Scalar, Vector}
 - Broadcast
 - Barrier

Brief description of MPI

- Point-to-point communications
 - For nearest neighbor, the requirements for QMP are well satisfied by QCDOC physics network
 - Both C/C++ Binдинг implemented over QCDOC system calls.
 - Store-and-forward capability of QCDOC is used.
 - Both C/C++ Binдинг implemented over QCDOC system calls.
- Global operations
 - Small software overhead.
 - Both C/C++ Binдинг implemented over QCDOC system calls.
 - Both C/C++ Binдинг implemented over QCDOC system calls.
 - Both C/C++ Binдинг implemented over QCDOC system calls.
- Implementation complete except non-nearest communication.
 - $L = \text{Number of sites } N^d = \text{Number of Dimension}$
 - Communication time grows as $\sim L \times N^d$
 - Both C/C++ Binдинг implemented over QCDOC system calls.
 - Implementation complete except non-nearest communication.

Implementation of QMP on QCDOC

While repeated communications are often to/from nearest-neighbor, NNC gives more flexibility and portability.

Non-Nearest Neighbor Communication(NNC) on QCDOC

- Non-interrupting physics network used for user data transfer
 - Interrupting communication channel (Supervisor) used for the header transfer
- Implementation strategy:
- Consideration for MPI implementation (R. Bennett, BNL).
- software overhead to nearest communication.
- Should be able to co-exist with nearest neighbor communication without introducing excessive overhead when destination is nearest
 - Minimum number of copying

Desired features:

- Consideration for MPI: Many of the necessary functionalities needed for MPI such as
 - Static, Manhattan-style routing to prevent lock-ups/long delays
 - Consideration for MPI: Many of the necessary functionalities needed for MPI such as
 - Message routing
 - Buffer management
 - Packetization
 - Queuing

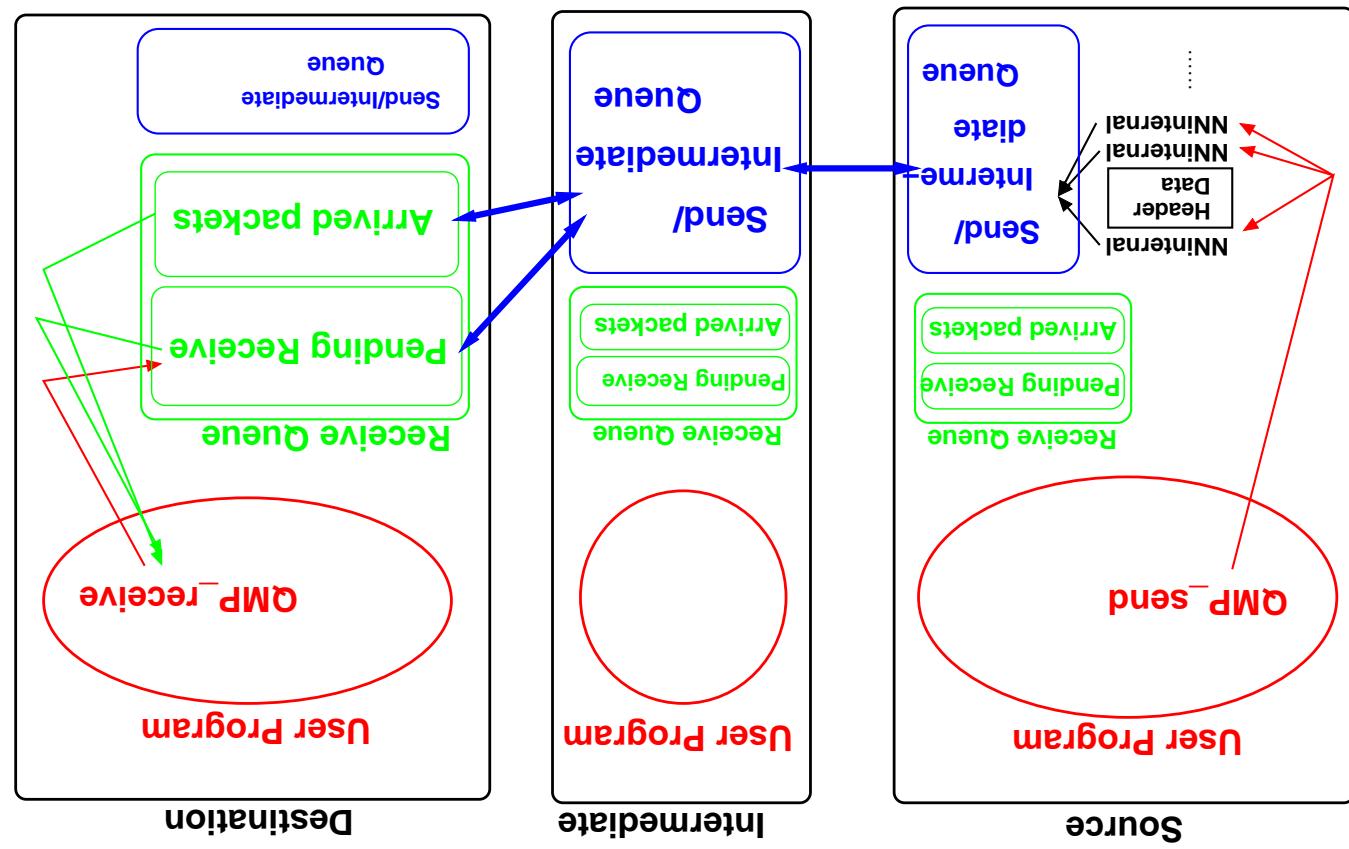


Diagram of NNC

- Linpack(D. Krohn, P. Boyle), etc..
- Runtime environment
- Parallel file system
- Batch system(PBS)
- Level 3 interfaces
- Implementation & Revision of MPI

Planned activities